# Recommendation system:
# Connecting business users with innovative solutions

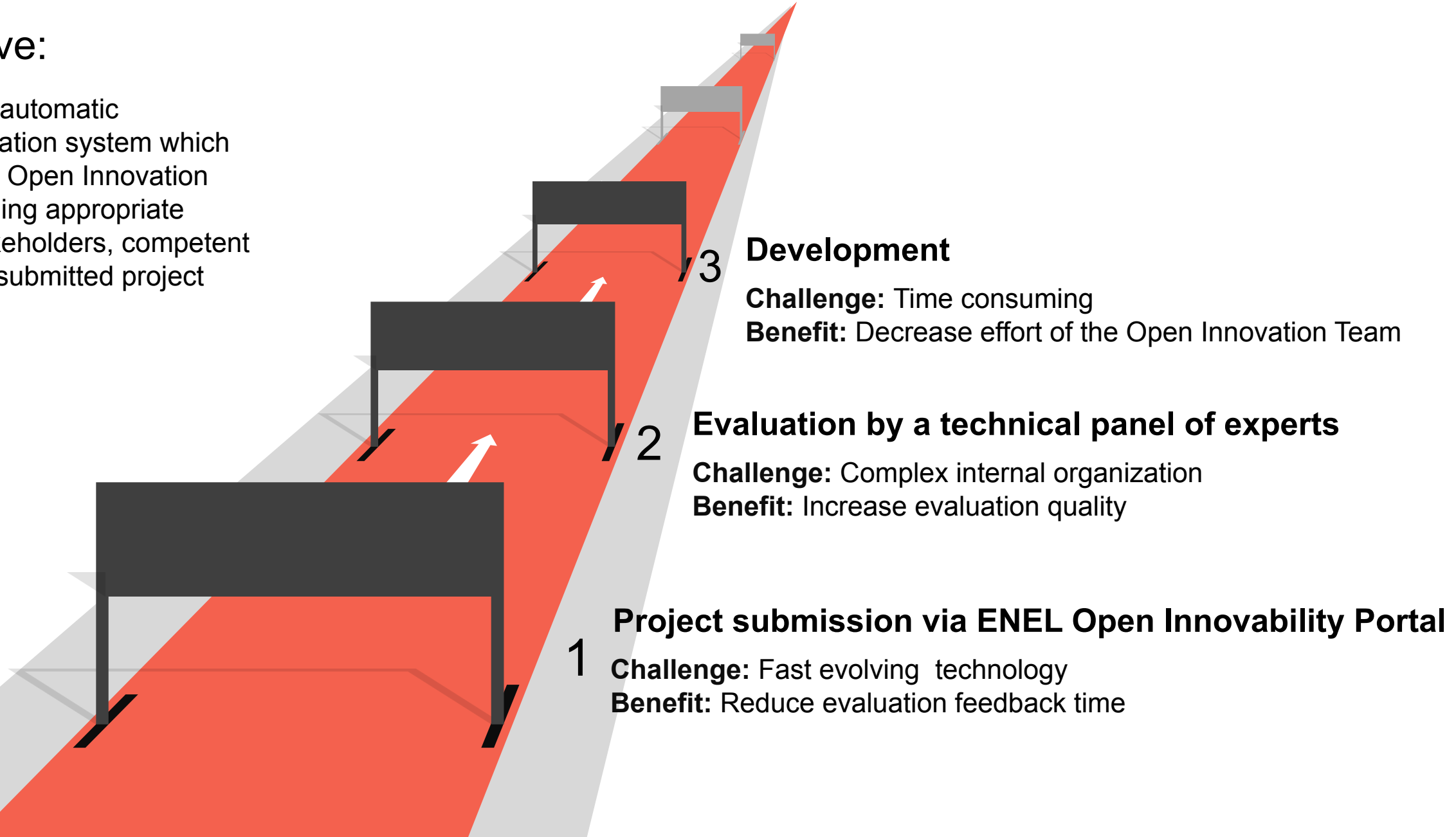**Candidate: Ivana Nastasic**

**Master Course: Data Science**

**Thesis Advisor: Prof. Pierpaolo Brutti**
**External Advisor: Marco Piersanti**

**Sapienza università di Roma - Facoltà di ingegneria dell'informazione informatica e statistica**

# Introduction

## Objective:

Develop an automatic recommendation system which will assist to Open Innovation Team in finding appropriate internal stakeholders, competent to evaluate submitted project proposal.



**3 Development**

**Challenge:** Time consuming
**Benefit:** Decrease effort of the Open Innovation Team

**2 Evaluation by a technical panel of experts**

**Challenge:** Complex internal organization
**Benefit:** Increase evaluation quality

**1 Project submission via ENEL Open Innovability Portal**

**Challenge:** Fast evolving technology
**Benefit:** Reduce evaluation feedback time

# Datasets

Row data export from The Open Innovability Portal, ~3800 Project Descriptions in several languages.

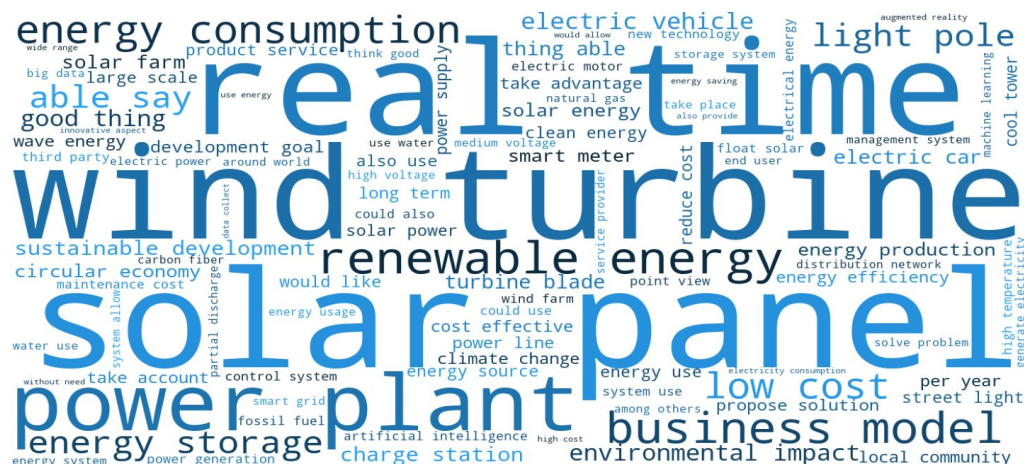Project Proposals dataset example

| | |
|---|---|
| SOL-27358 | There is a plugin Office called "Dictate" . It can be downloaded from this Microsoft website ( dictate.ms ) . Using this plugin Office programs (Outlook,Word,Powerpoint...) can write automatically or translate. in another Language. The benefit is that for an Enel employee, is easier and faster, to think a document and to speak to this "digital secretary", then to type on the desk. To be more clearer, please look at this 2 youtube walk-through videos: https://www.youtube.com/watch?v=auF9bvAectU https://www.youtube.com/watch?v=k9gCfEJGj38 |

Processed dataset based on self presentation of employees on internal e-profile portal, ~ 35000 employees, 18 skill types and 298 skill subtypes. Skills are entered in free text format.
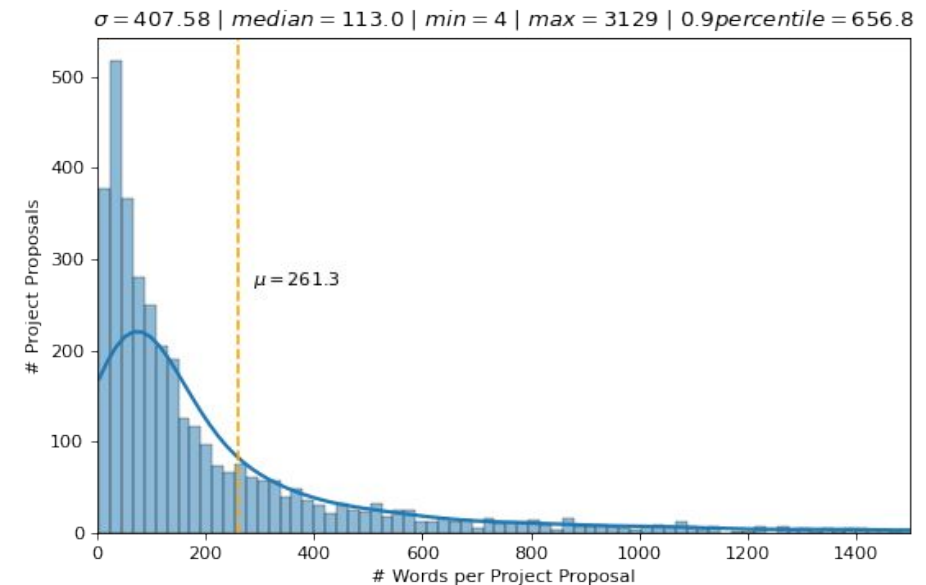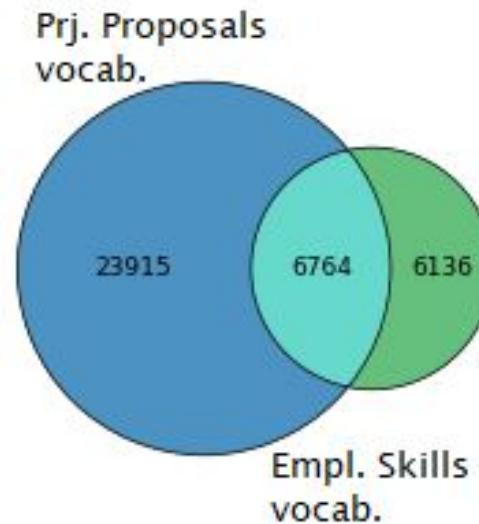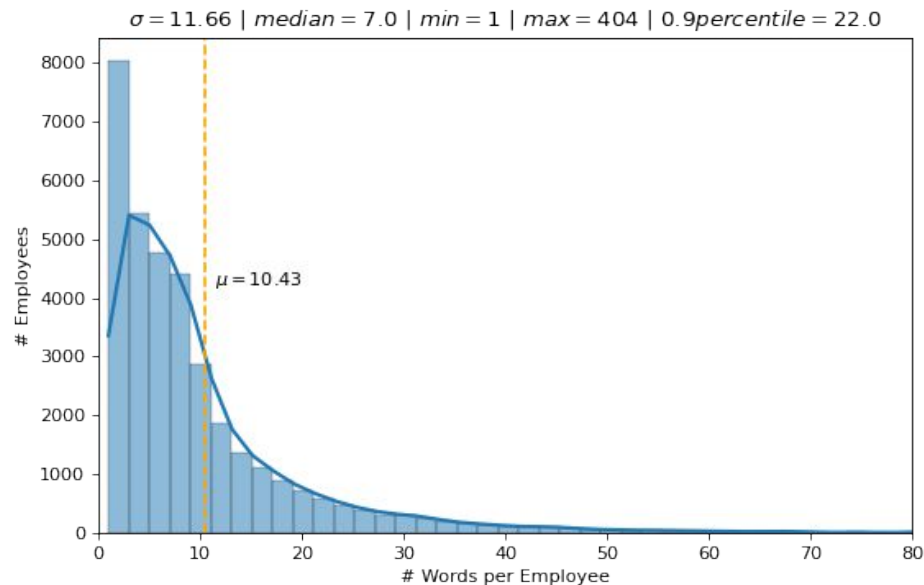
Employee's skills dataset example

| | |
|---|---|
| Employee ID:196 Skill ID: 1131248816-2 | Skill type: energy related skills Skill subtype: power plants Skill description: power generation management |

# Exploratory Data Analysis Findings

**01** Text descriptions of employee's skills are very short.

**02** Vocabularies are differ significantly in size and content.

**03** There is a huge difference in text length between projects and employees descriptions.



$\sigma = 11.66 \mid median = 7.0 \mid min = 1 \mid max = 404 \mid 0.9 percentile = 22.0$

$\mu = 10.43$

# Words per Employee

Prj. Proposals vocab.

23915    6764    6136

Empl. Skills vocab.

$\sigma = 407.58 \mid median = 113.0 \mid min = 4 \mid max = 3129 \mid 0.9 percentile = 656.8$

$\mu = 261.3$

# Words per Project Proposal

# Evaluation

**01** User-centric Perceived Recommendation Accuracy: % of project proposals with at least 1 relevant suggestion.

**20** Project Proposals
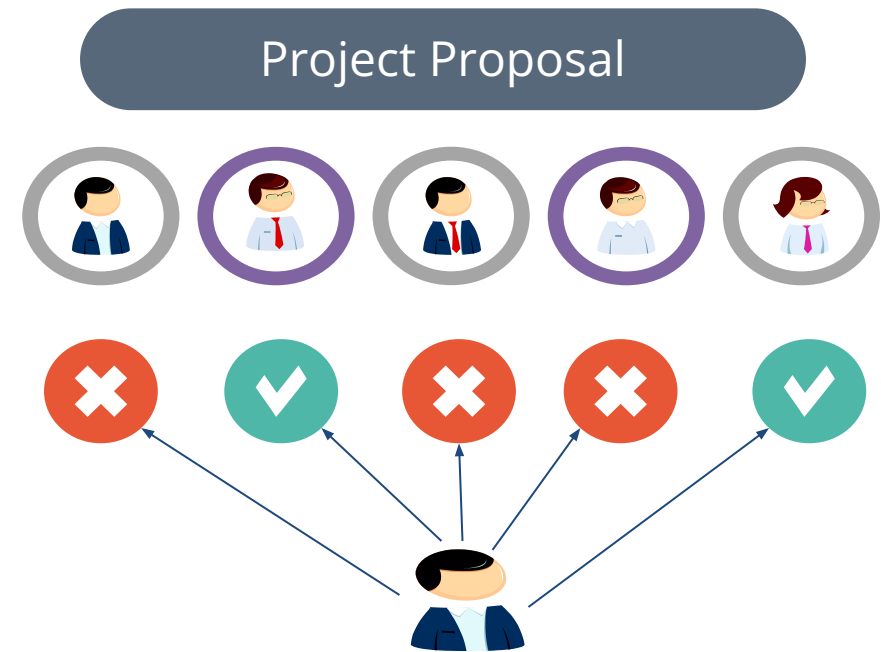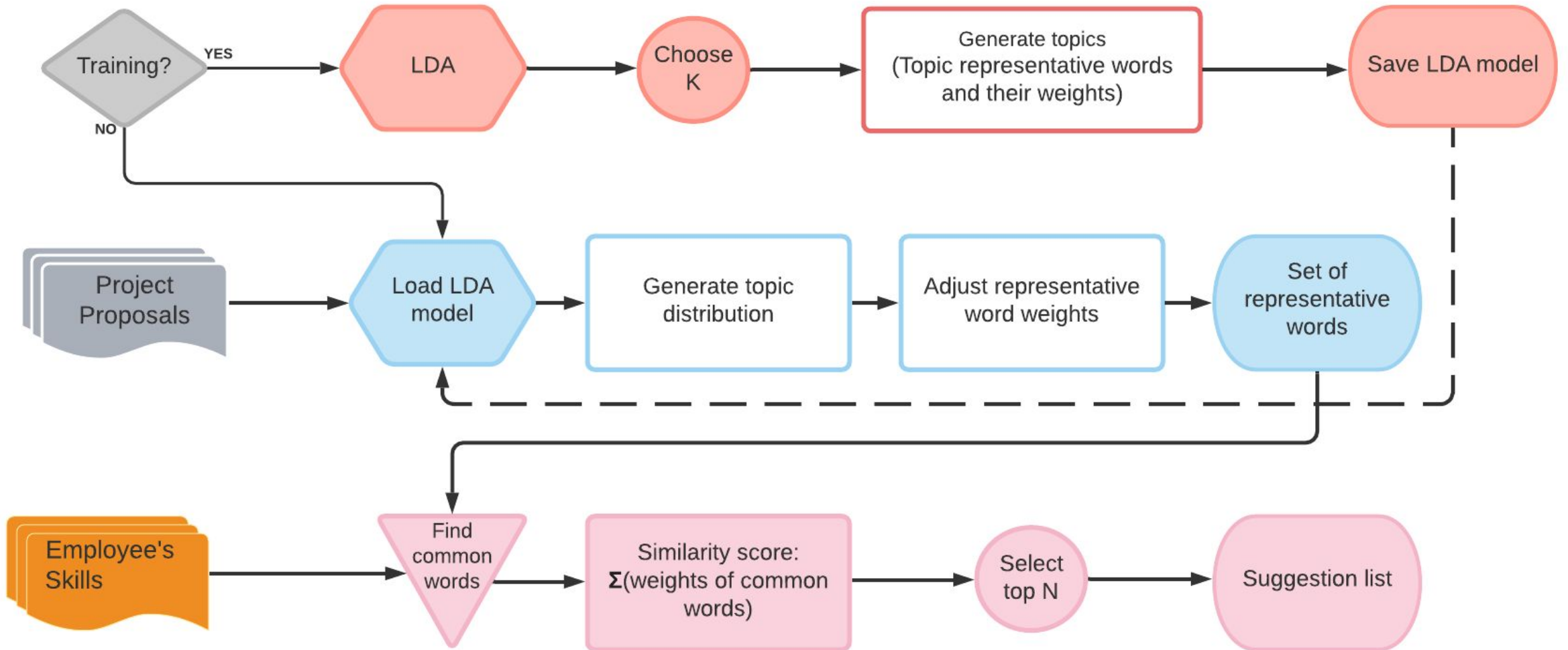
**5** Suggested Employees

**10** Perceived Relevance
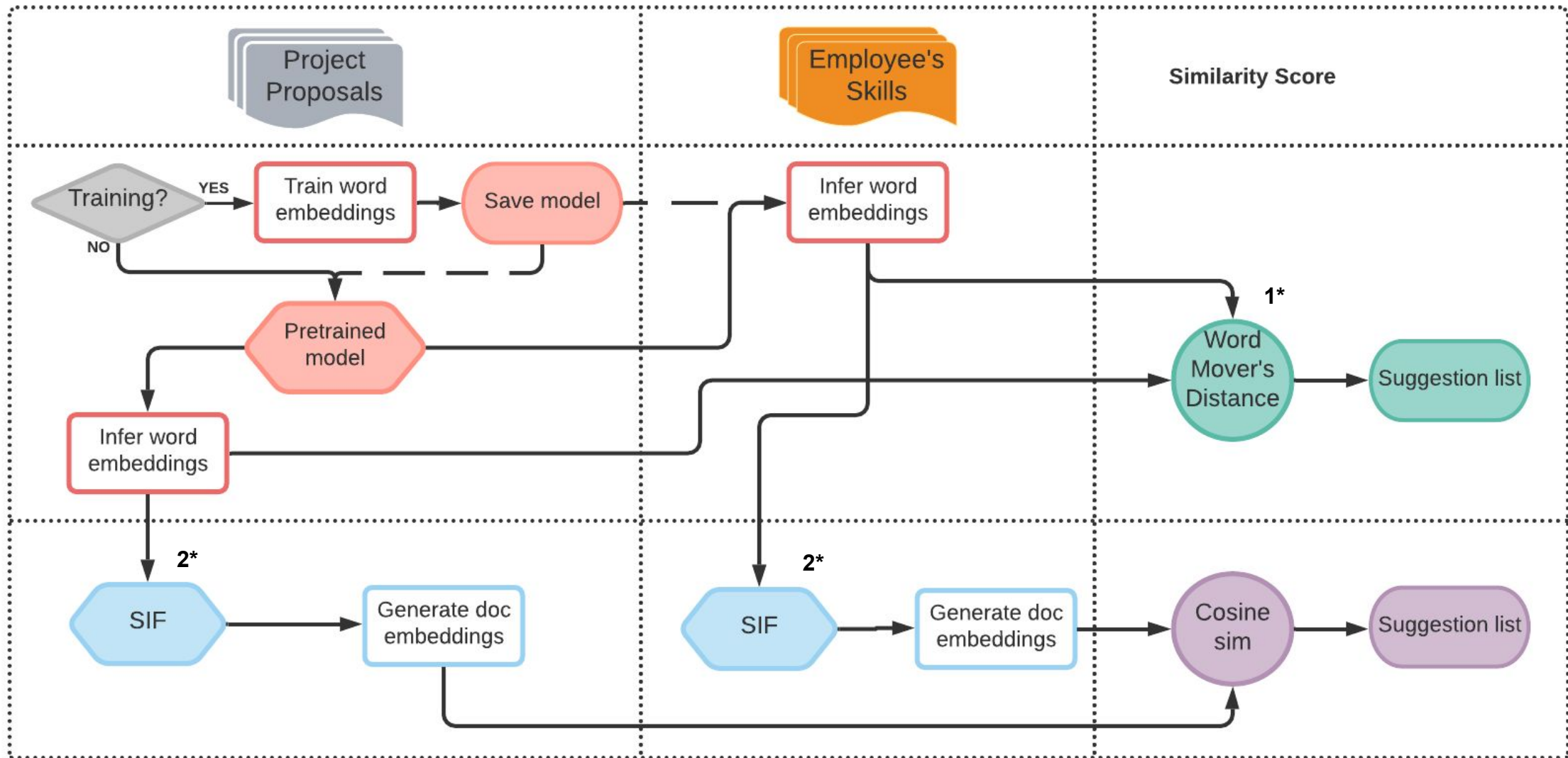
Open Innovation Team

**3** Model setups

Project Proposal

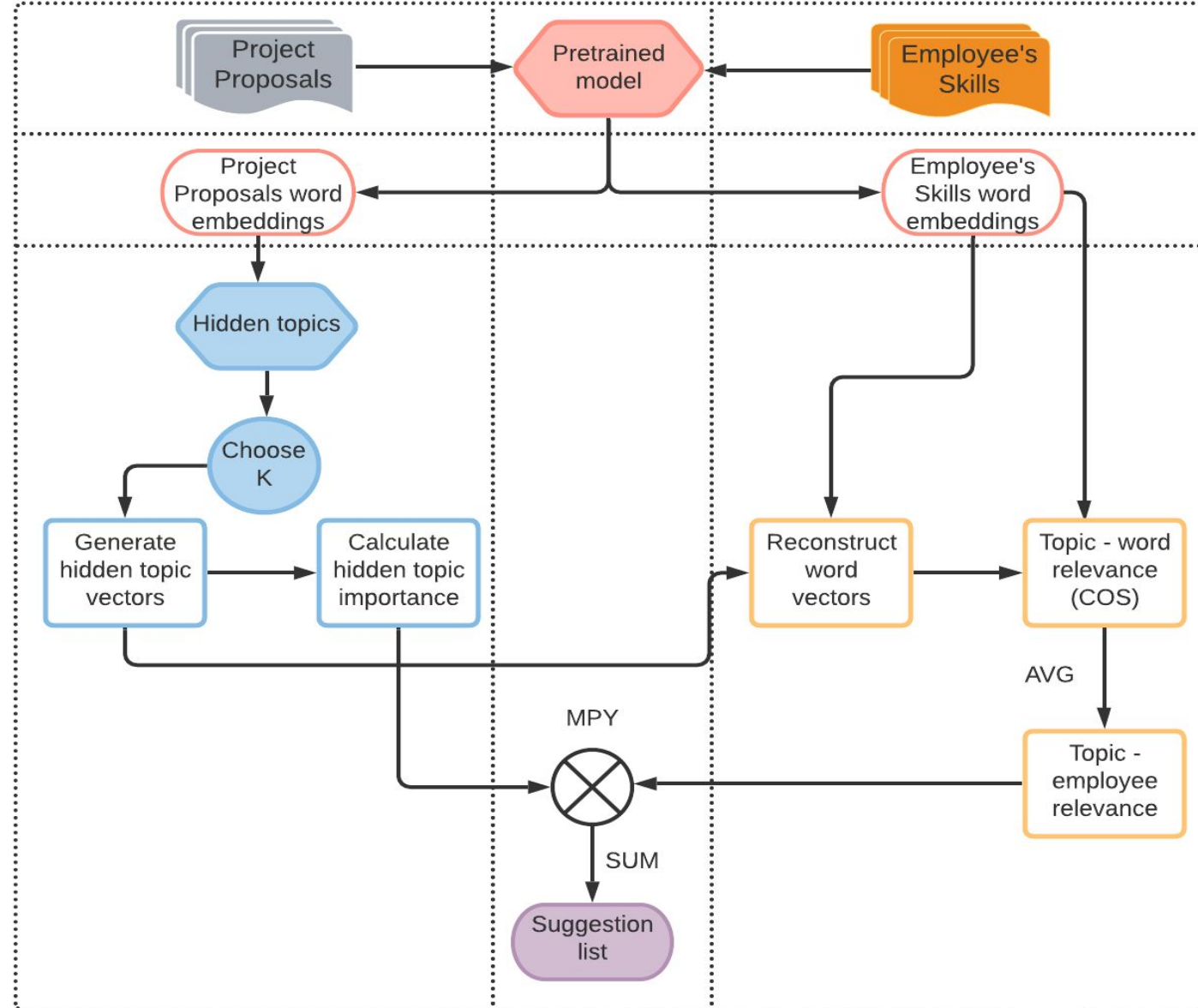# LDA topic modeling and significant words matching

# Text embeddings for similarity calculation



Papers: 1*.  Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. From word embeddings to document distances, 2015

2*. Sanjeev Arora, Yingyu Liang, and Tengyu Ma. A simple but tough-to-beat baseline for sentence embeddings, 2017

# Matching texts of varying length via hidden topics



**Paper: Hongyu Gong, Tarek Sakakini, Suma Bhat, and Jinjun Xiong. Document similarity for texts of varying lengths via hidden topics, 2019**

# Results

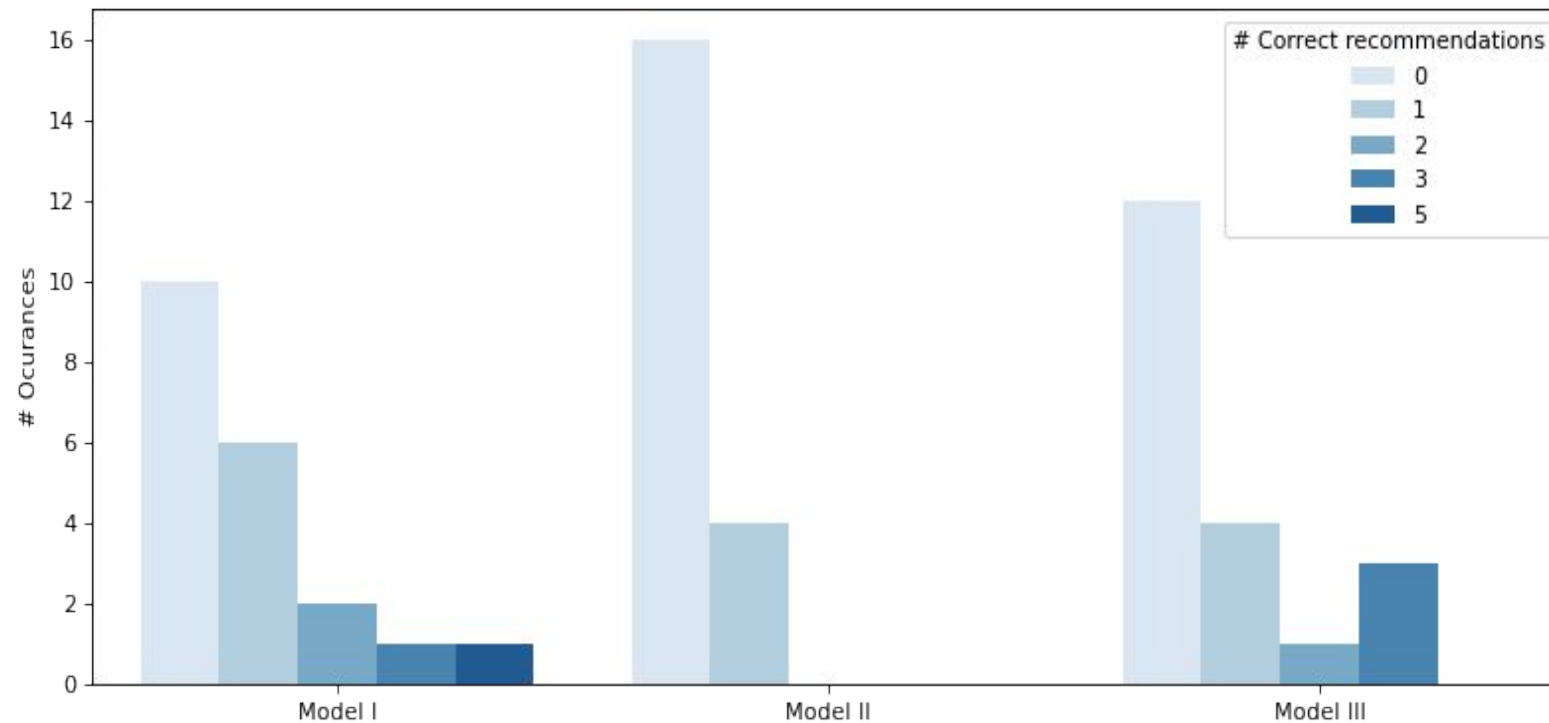| | | |
|---|---|---|
| **0.5** | **Model I**<br>LDA topic modeling and significant words matching with number of topics K=90. | |
| **0.25** | **Model II**<br>SIF with self trained word embeddings using fastText model, vector dimension d=300. | |
| **0.4** | **Model III**<br>Hidden topics with pretrained word embeddings trained with fastText model on Common Crawl dataset* , vector dim d=300 and  K=5. | |



*https://fasttext.cc/docs/en/english-vectors.html

# Conclusions

**01**    **Text embeddings performed poorly due to the big difference in text lengths.**

**02**    **Hidden topics approach needs tuning of number of topics K.**

**03**    **The best scored method is: LDA topic modeling and significant words matching.**
**Further development:**
- **Full submitted project documentation**
- **Standardized employee's skills (ESCO, O*NET)**

Thank you!